

- Ill conditioned Problem (multicollinearity ($X^T X$) is singular) : $\hat{\beta}_{OLS} = (X^T X)^{-1} X^T Y$
- To fix it we use Ridge, LASSO or Elastic Net

* **Ridge Regression Method :**

$$L^2 \text{ loss} = \sum_{i=1}^n e_i^2 = SSE$$

$$L^1 \text{ loss} = \sum_{i=1}^n |e_i| \quad \text{Not differentiable}$$

$$\hat{\beta}_{OLS} = \min_{\beta} (L^2 \text{ loss}) * \text{ see we get } \hat{\beta}_{OLS}$$

Penalty terms ensures that $(X^T X)^{-1}$ always exists

* is reformulated to take loss + penalty approach that attends to shrink the coefficients (some of the parameters are close to zero) by imposing a penalty on their size (size of the coefficients).

$$Y = \beta_0 + \sum_{i=1}^p \beta_i X_i + \varepsilon$$

$$L^2 \text{ loss} = \sum_{i=1}^n (Y_i - \beta_0 - \sum_{j=1}^p \beta_j X_j)^2$$

Two formulations were proposed :

$$\hat{\beta}_{\text{Ridge}} = \min_{\beta} \left\{ \sum_{i=1}^n (\gamma_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij})^2 + \lambda \sum_{j=1}^p \beta_j^2 \right\}$$

$$\hat{\beta}_{\text{Ridge}} = \min_{\beta} \left\{ \sum_{i=1}^n (\gamma_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij})^2 \right\}; \sum_{j=1}^p \beta_j^2 \leq t, t \neq 1$$

λ := shrinkage parameter estimated by cross validation.

$$\lambda = 0 \Rightarrow \hat{\beta}_{\text{Ridge}} = \hat{\beta}_{\text{OLS}}$$

Reel absent Time varying Ridge Model

* Estimationen at $\hat{\beta}$:

$$\min ((Y - X\beta)^T (Y - X\beta) + \lambda \beta^T \beta) =$$

$$\min L(Y^T - \beta^T X^T)(Y - X\beta) + \lambda \beta^T \beta =$$

$$\min (Y^T Y - Y^T X \beta - \beta^T X^T Y + \beta^T X^T X \beta + \lambda \beta^T \beta) =$$

$$\min (Y^T Y - \beta^T X^T Y - (\beta^T X^T Y)^T + \beta^T X^T X \beta + \lambda \beta^T \beta) =$$

$$\min (Y^T Y - 2\beta^T X^T Y + \beta^T X^T X \beta + \lambda \beta^T \beta)$$

$$\frac{\partial}{\partial \beta} = -2X^T Y + 2X^T X \beta + 2\lambda \beta = 0$$

$$X^T X \beta + \lambda \beta - X^T Y = 0$$

$$\beta (X^T X + \lambda I) = X^T Y$$

$$\hat{\beta}_{\text{Ridge}} = (X^T X + \lambda I)^{-1} X^T Y$$

- * The coefficients can be estimated by OLS & λ is estimated by cross validation.
- * The choice of quadratic penalty adds a term proportional to the diagonal element of $(X^T X)$ and that mitigates the problem of multicollinearity.
- * $\hat{\beta}_{\text{Ridge}}$ is linear

$$\begin{aligned}
 * E(\hat{\beta}_{\text{Ridge}}) &= E((X^T X + \lambda I)^{-1} X^T Y) \\
 &= (X^T X + \lambda I)^{-1} X^T E(X\beta + \varepsilon) \\
 &= (X^T X + \lambda I)^{-1} X^T X \beta \\
 &\neq \beta
 \end{aligned}$$

Hence, it is a biased estimator but in the case of bias, ridge reduces the variance

$$\begin{aligned}
 \text{Var}(\hat{\beta}_{\text{Ridge}}) &= \text{Var}((X^T X + \lambda I)^{-1} X^T Y) \\
 &= ((X^T X + \lambda I)^{-1} X^T)((X^T X + \lambda I)^{-1} X^T)^T \text{Var}(Y) \\
 &= \sigma^2 (X^T X + \lambda I)^{-1} X^T X (X^T X + \lambda I)^{-1}
 \end{aligned}$$

* LASSO Regression :

LASSO : Least absolute shrinkage & selection operator

In Ridge the penalty is $\lambda \sum_{j=1}^p \beta_j^2$ shrinks all $\beta_j \rightarrow 0$ but LASSO solves this problem.

LASSO uses L₁ norm that has the effect of forcing some coefficients to become zero which leads to sparse model.

$$\min_{\beta} \left(\sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij})^2 + \lambda \sum_{j=1}^p |\beta_j| \right)$$

$$\min_{\beta} \left(\sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij})^2 \right); \sum_{j=1}^p |\beta_j| \leq t, t > 0$$

λ is estimated by cross validation.

$\hat{\beta}_{\text{LASSO}}$ does not have a closed solution like OLS & Ridge but it is helpful when we have 1000 features but only 200 samples as an example.

We can use package to do this year!

Also, check Bayesian LASSO, Polynomial, Non linear and L₀ regularization

* Elastic Net Method :

$$\hat{\beta}_{\text{EN}} = \text{Loss} + \lambda \sum_{j=1}^p (\alpha \beta_j^2 + (1-\alpha) |\beta_j|)$$

so it is a compromise between Ridge & LASSO

Ridge

LASSO

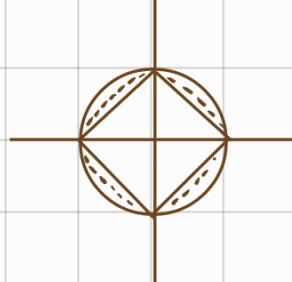
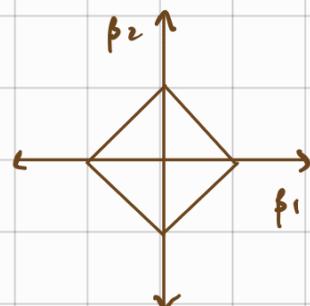
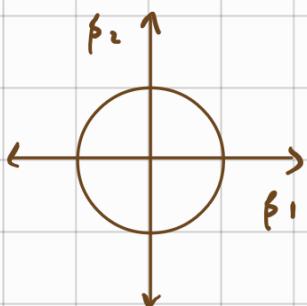
Elastic Net

$$\beta_1^2 + \beta_2^2 \leq 1$$

$$u^2 + y^2 \leq 1$$

$$|\beta_1| + |\beta_2| \leq 1$$

$$|u| + |y| \leq 1$$



* Remark : Before fitting regularization models (LASSO, Ridge) the predictors x 's should be standanized (unless using a package)

$x_{\text{new}} = \frac{x - \text{mean}}{\text{sd}}$ because LSE is a scaling variant but this is not.

glmnet standardize by default

* Standard error : \sqrt{MSRes}

* We use cross validation to estimate alpha and lambda for elastic net in R for best results (10 fold CV).

* Transformations :

| | x | y | \hat{y} |
|-------|-------------------------|-----------------------|---------------|
| Train | $\log(x_{\text{Tr}})$ | $\log(y_{\text{Tr}})$ | |
| Test | $\log(x_{\text{Test}})$ | \hat{y} | $e^{\hat{y}}$ |

Reverse transform

* Standardization Comparison : $\frac{n - \mu_n}{\sigma_n}$

Train : $\frac{y_{Tr} - \mu_{yTr}}{\sigma_{yTr}}$

Test : $\frac{y_{Ts} - \mu_{yTs}}{\sigma_{yTs}}$

don't use μ_{Ts} / σ_{Ts}
otherwise we
get data leakage

High Range / Not following assumptions :

$$y' \longrightarrow \frac{y - \mu_y}{\sigma_y}$$

$$y'_{Tr} = \frac{y_{Tr} - \mu_{yTr}}{\sigma_{yTr}}$$

$$y'_{Ts} = \frac{y_{Ts} - \mu_{yTs}}{\sigma_{yTs}}$$